

Sept 20<sup>th</sup>

Feature Space - Set of all possible data points.

COMPARISON REQUIRES DISSIMILARITY FUNC [input 2 data points, large value  $\rightarrow$  v. dissimilar]

\*one instance is a distance func.

MINKOWSKI.

ALSO CAN USE SIMILARITY FUNC [large value  $\rightarrow$  v. similar]

COSINE SIMILARITY

$$S(x, y) = \cos \theta \quad \text{where } \theta \text{ is angle b/w } x, y$$

WE CAN GET CORRESPONDING DISSIM FUNC

$$d(x, y) = \frac{1}{S(x, y)} \quad \text{OR} \quad d(x, y) = k - S(x, y) \quad \text{for some } k.$$

$$d(x, y) = 1 - S(x, y)$$

WE USE COSINE OVER EUCLIDEAN DISTANCE?  $\rightarrow$  WHEN DIRECTION MATTERS MORE THAN MAGNITUDE

Jaccard Similarity

ex: take text and each word is an attribute, and values are whether it is in the given doc.

$\rightarrow$  Could use MANHATTAN  $\rightarrow$  RETURNS SIZE OF SET DIFFERENCE.

PROBLEM IS IT CAN'T DISTINGUISH HOW MANY HAD THE SAME VALUE.

WE WANT SOME NOTION OF INTERSECTION

$$Jsim(x, y) = \frac{|x \cap y|}{|x \cup y|} \quad JDist(x, y) = 1 - \frac{|x \cap y|}{|x \cup y|}$$